



# Chapter 3

# Dynamic Programming

## Part 2

Algorithm Theory  
WS 2012/13

Fabian Kuhn

# Dynamic Programming

---



„Memoization“ for increasing the efficiency of a recursive solution:

- Only the *first time* a sub-problem is encountered, its **solution is computed** and then stored in a table. Each subsequent time that the subproblem is encountered, the value stored in the table is simply looked up and returned  
(without repeated computation!).
- **Computing the solution**: For each sub-problem, store how the value is obtained (according to which recursive rule).

# Dynamic Programming

---

Dynamic programming / memoization can be applied if

- **Optimal solution** contains **optimal solutions to sub-problems**  
(recursive structure)

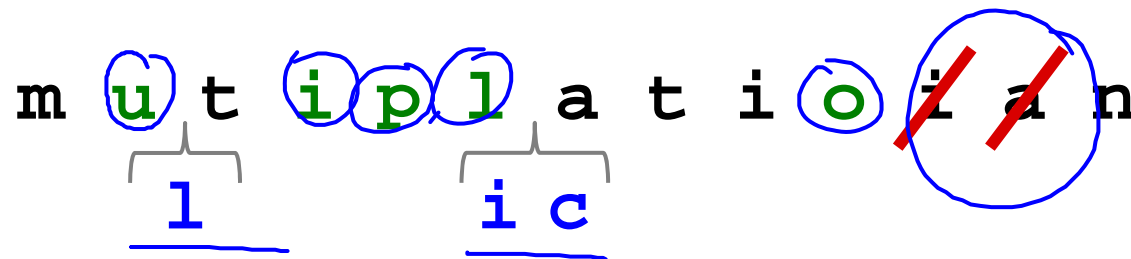
- **Number** of sub-problems that need to be considered is small  
*polynomial*

# String Matching Problems

## Edit distance:

- For two given strings  $A$  and  $B$  efficiently compute the **edit distance**  $D(A, B)$  (# edit operations to transform  $A$  into  $B$ ) as well as a minimum sequence of edit operations that transform  $A$  into  $B$ .

- Example:** mathematician  $\rightarrow$  multiplication:



# String Matching Problems

**Edit distance  $D(A, B)$**  (between strings  $A$  and  $B$ ):

m a - t h e m - - a t i c i a n  
m u l t i p l i c a t i o - - n

**Approximate string matching:**

For a given text  $T$ , a pattern  $P$  and a distance  $d$ , find all  
substrings  $P'$  of  $T$  with  $D(P, P') \leq d$ .

**Sequence alignment:**

Find optimal alignments of DNA / RNA / ... sequences.

G A G C A - C T T G G A T T C T C G G  
- - - C A C G T G G - A - A C T - - -

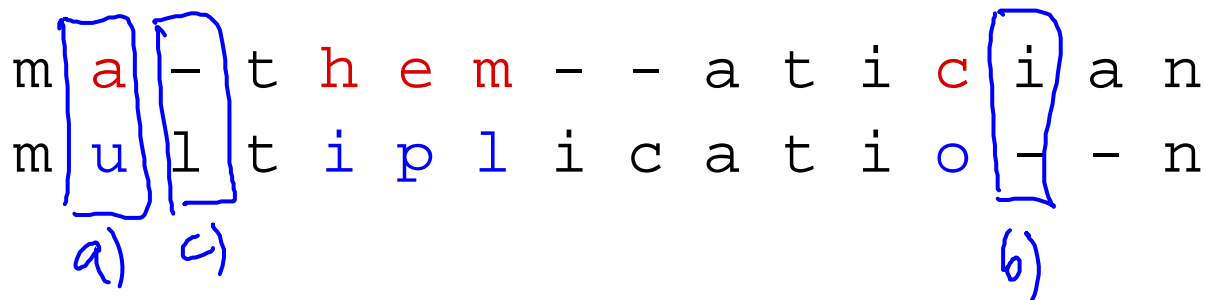
# Edit Distance

**Given:** Two strings  $A = a_1 a_2 \dots a_m$  and  $B = b_1 b_2 \dots b_n$

**Goal:** Determine the minimum number  $D(A, B)$  of edit operations required to transform  $A$  into  $B$

## Edit operations:

- a) **Replace** a character from string  $A$  by a character from  $B$
- b) **Delete** a character from string  $A$
- c) **Insert** a character from string  $B$  into  $A$



# Edit Distance – Cost Model

- Cost for **replacing** character  $a$  by  $b$ :  $c(a, b) \geq 0$
- Capture insert, delete by allowing  $a = \underline{\varepsilon}$  or  $b = \underline{\varepsilon}$ :
  - Cost for **deleting** character  $a$ :  $c(a, \varepsilon)$
  - Cost for **inserting** character  $b$ :  $c(\varepsilon, b)$

- **Triangle inequality:**

$$\underline{c(a, c) \leq c(a, b) + c(b, c)}$$

→ each character is changed at most once!

- **Unit cost model:**  $c(a, b) = \begin{cases} 1, & \text{if } a \neq b \\ 0, & \text{if } a = b \end{cases}$

# Recursive Structure

- Optimal “alignment” of strings (unit cost model)

bbcadfagikccm and abbagflrgikacc:

A	-	b	b	c	a	g	f	a	-	g	i	k	-	c	c	m
B	a	b	b	-	a	d	f	l	r	g	i	k	a	c	c	-
	ms.			del		repl.		repl.	ins				ins			del

- Consists of optimal “alignments” of sub-strings, e.g.:

-bbcagfa      and      -gik-ccm  
 abb-adfl      rgikacc-

$$A_{ij} = a_i \dots a_j$$

- Edit distance between  $A_{1,m} = a_1 \dots a_m$  and  $B_{1,n} = b_1 \dots b_n$ :

$$D(A, B) = \min_{\substack{k, \ell \\ k < m, \ell < n}} \{ \underline{D(A_{1,k}, B_{1,\ell})} + \underline{D(A_{k+1,m}, B_{\ell+1,n})} \}$$

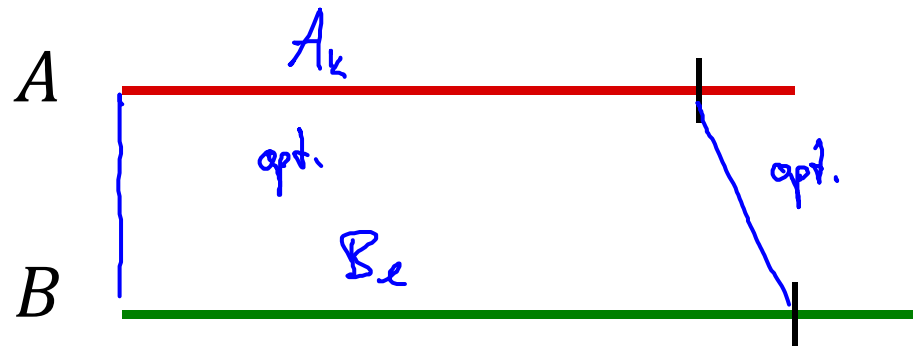


# Computation of the Edit Distance

Let  $A_k := a_1 \dots a_k$ ,  $B_\ell := b_1 \dots b_\ell$ , and

$$D_{k,\ell} := D(A_k, B_\ell)$$

$$D_{m,n} = D(A, B)$$

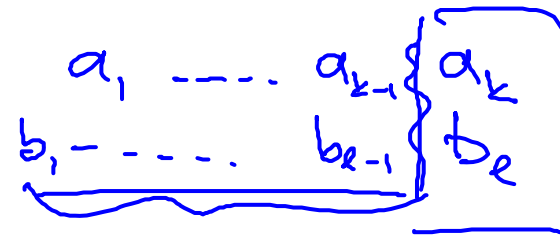


# Computation of the Edit Distance

Three ways of ending an “alignment” between  $A_k$  and  $B_\ell$ :

1.  $a_k$  is replaced by  $b_\ell$ :

$$\underline{D_{k,\ell}} = \underline{D_{k-1,\ell-1}} + \underline{c(a_k, b_\ell)}$$



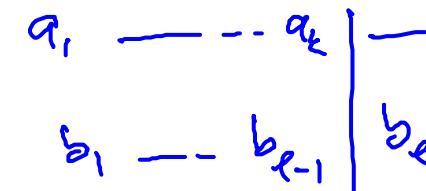
2.  $a_k$  is deleted:

$$\underline{D_{k,\ell}} = \underline{D_{k-1,\ell}} + \underline{c(a_k, \varepsilon)}$$



3.  $b_\ell$  is inserted:


$$\underline{D_{k,\ell}} = \underline{D_{k,\ell-1}} + \underline{c(\varepsilon, b_\ell)}$$



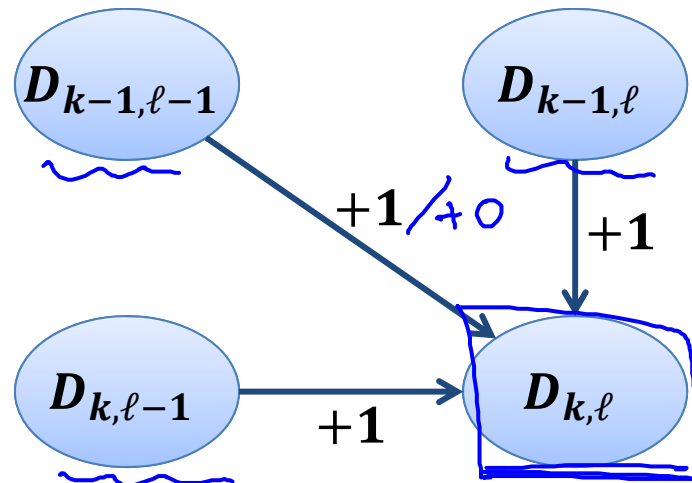
# Computing the Edit Distance

- Recurrence relation (for  $k, \ell \geq 1$ )

$$\underline{D_{k,\ell}} = \min \left\{ \begin{array}{l} D_{k-1,\ell-1} + c(a_k, b_\ell) \\ D_{k-1,\ell} + c(a_k, \varepsilon) \\ D_{k,\ell-1} + c(\varepsilon, b_\ell) \end{array} \right\} = \min \underbrace{\left\{ \begin{array}{l} D_{k-1,\ell-1} + \mathbf{1} \\ D_{k-1,\ell} + 1 \\ D_{k,\ell-1} + 1 \end{array} \right\}}_{\text{unit cost model}}$$

*0 if  $a_k = b_\ell$*   


- Need to compute  $D_{i,j}$  for all  $0 \leq i \leq k, 0 \leq j \leq \ell$ :



# Recurrence Relation for the Edit Distance



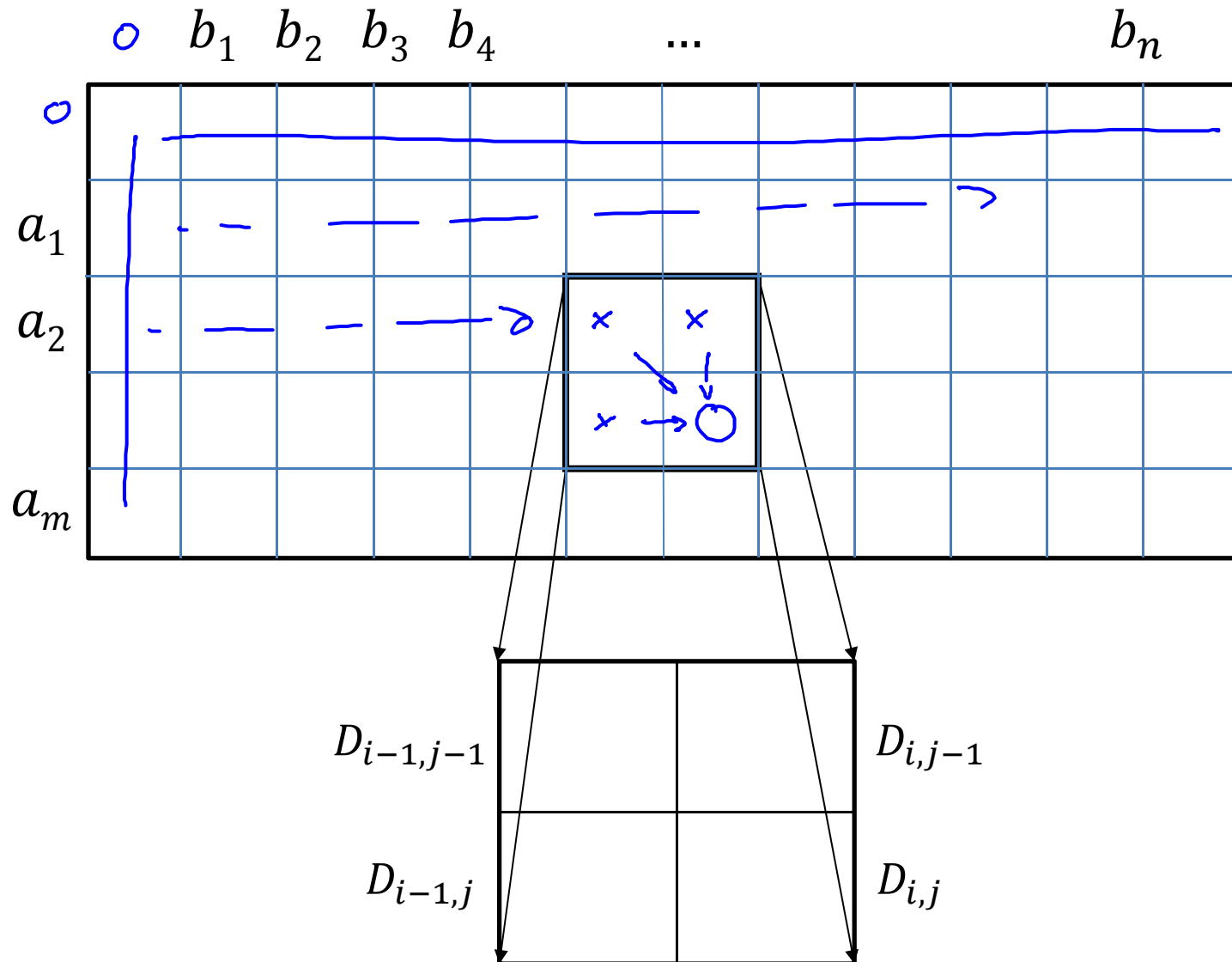
Base cases:

$$\begin{aligned} D_{0,0} &= D(\varepsilon, \varepsilon) = 0 && \text{unit cost} \\ D_{\underline{0},j} &= D(\varepsilon, B_j) = D_{0,j-1} + \underline{c(\varepsilon, b_j)} = j \\ D_{i,\underline{0}} &= D(A_i, \varepsilon) = D_{i-1,0} + c(a_i, \varepsilon) = i \end{aligned}$$

Recurrence relation:

$$D_{i,j} = \min \left\{ \begin{array}{l} D_{k-1,\ell-1} + c(a_k, b_\ell) \\ D_{k-1,\ell} + c(a_k, \varepsilon) \\ D_{k,\ell-1} + c(\varepsilon, b_\ell) \end{array} \right\}$$

# Order of solving the subproblems



# Algorithm for Computing the Edit Distance



## Algorithm *Edit-Distance*

**Input:** 2 strings  $A = a_1 \dots a_m$  and  $B = b_1 \dots b_n$

**Output:** matrix  $D = (D_{ij})$

1  $D[0,0] := 0;$

2 **for**  $i := 1$  **to**  $m$  **do**  $D[i, 0] := i;$

3 **for**  $j := 1$  **to**  $n$  **do**  $D[0, j] := j;$

4 **for**  $i := 1$  **to**  $m$  **do**  $\leftarrow$  rows

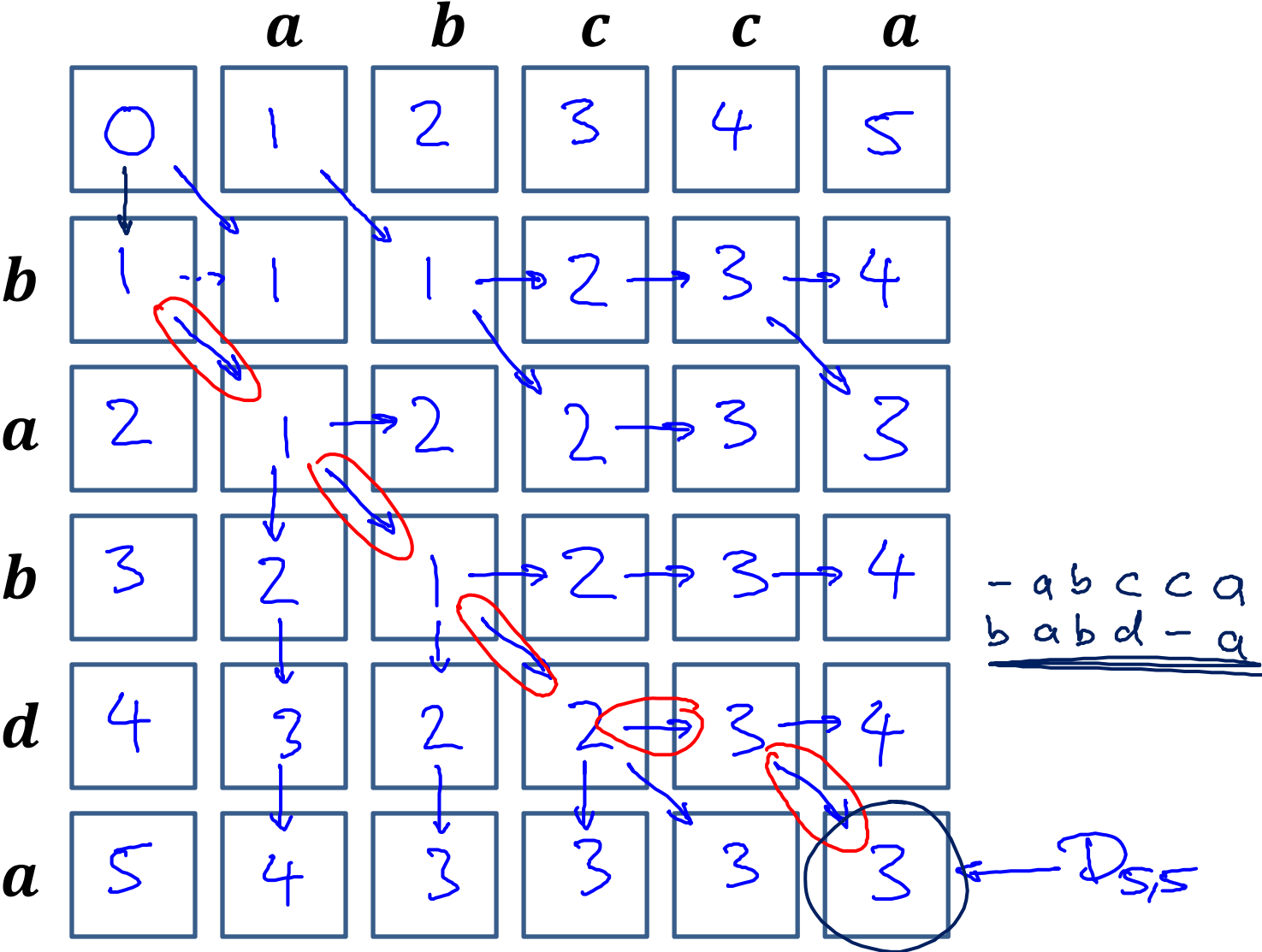
5 **for**  $j := 1$  **to**  $n$  **do**  $\leftarrow$  columns

6  $D[i, j] := \min \left\{ \begin{array}{l} D[i-1, j] + 1 \\ D[i, j-1] + 1 \\ D[i-1, j-1] + c(a_i, b_j) \end{array} \right\};$

} initialize

# Example

unit cost



# Computing the Edit Operations

**Algorithm** *Edit-Operations*( $i, j$ )

**Input:** matrix  $D$  (already computed)

**Output:** list of edit operations

- 1 **if**  $i = 0$  **and**  $j = 0$  **then return** empty list
- 2 **if**  $i \neq 0$  **and**  $D[i, j] = \underline{D[i - 1, j] + 1}$  **then**
- 3     **return** *Edit-Operations*( $i - 1, j$ )  $\circ$  „delete  $a_i$ “
- 4 **else if**  $j \neq 0$  **and**  $D[i, j] = D[i, j - 1] + 1$  **then**
- 5     **return** *Edit-Operations*( $i, j - 1$ )  $\circ$  „insert  $b_j$ “
- 6 **else** //  $D[i, j] = D[i - 1, j - 1] + c(a_i, b_j)$
- 7     **if**  $a_i = b_i$  **then return** *Edit-Operations*( $i - 1, j - 1$ )
- 8     **else return** *Edit-Operations*( $i - 1, j - 1$ )  $\circ$  „replace  $a_i$  by  $b_j$ “

**Initial call:** *Edit-Operations*( $m, n$ )



# Edit Operations

		<i>a</i>	<i>b</i>	<i>c</i>	<i>c</i>	<i>a</i>
	0	1	2	3	4	5
<i>b</i>	1	1	1	2	3	4
<i>a</i>	2	1	2	2	3	3
<i>b</i>	3	2	1	2	3	4
<i>d</i>	4	3	2	2	3	4
<i>a</i>	5	4	3	3	3	3

Red arrows indicate the path from (0,0) to (5,5): (0,0) → (1,0) → (1,1) → (2,1) → (2,2) → (3,2) → (3,3) → (4,3) → (4,4) → (5,4) → (5,5). The cell (5,5) containing the value 3 is circled in red.

# Edit Distance: Summary

- Edit distance between two strings of length  $m$  and  $n$  can be computed in  $O(mn)$  time.
- Obtain the edit operations:
  - for each cell, store which rule(s) apply to fill the cell
  - track path backwards from cell  $(m, n)$
  - can also be used to get all optimal “alignments”
- Unit cost model:
  - interesting special case
  - each edit operation costs 1

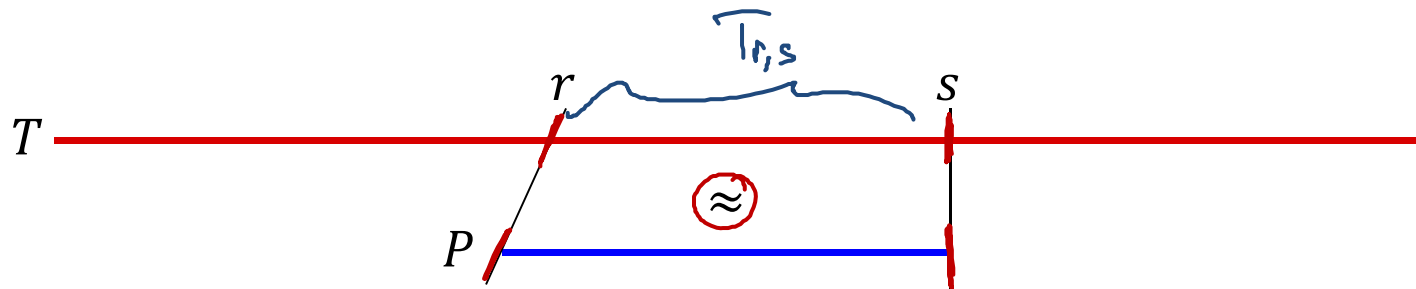
$O(m+n)$

# Approximate String Matching

**Given:** strings  $T = t_1 t_2 \dots t_n$  (text) and  $P = p_1 p_2 \dots p_m$  (pattern).

**Goal:** Find an interval  $[r, s]$ ,  $1 \leq r \leq s \leq n$  such that the sub-string  $T_{r,s} := t_r \dots t_s$  is the one with highest similarity to the pattern  $P$ :

$$\arg \min_{1 \leq r \leq s \leq n} D(T_{r,s}, P)$$



# Approximate String Matching



## Naive Solution:

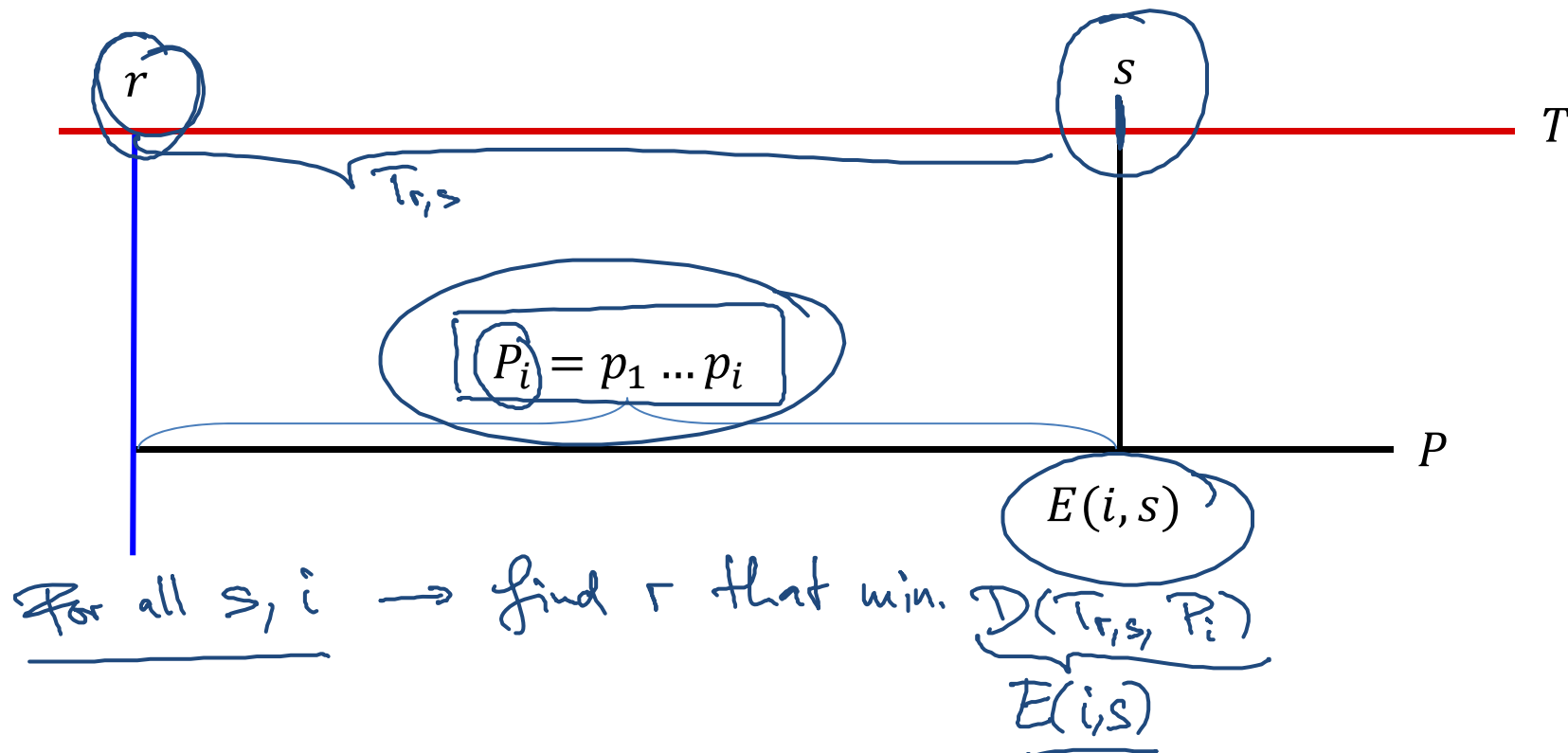
**for all**  $1 \leq r \leq s \leq n$  **do**  
    compute  $D(T_{r,s}, P)$   
choose the minimum

$$O(n^2 \cdot n \cdot m) = O(n^3 m)$$

# Approximate String Matching

A related problem:

- For each position  $s$  in the text and each position  $i$  in the pattern compute the minimum edit distance  $E(i, s)$  between  $P_i = p_1 \dots p_i$  and any substring  $T_{r,s}$  of  $T$  that ends at position  $s$ .



# Approximate String Matching

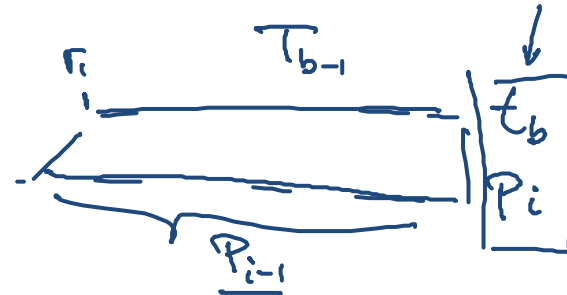
*S from previous slide*



Three ways of ending optimal alignment between  $T_b$  and  $P_i$ :

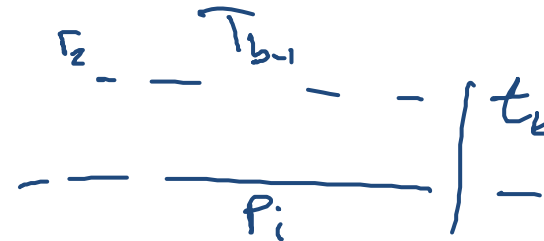
1.  $t_b$  is replaced by  $p_i$ :

$$E_{b,i} = E_{b-1,i-1} + c(t_b, p_i)$$



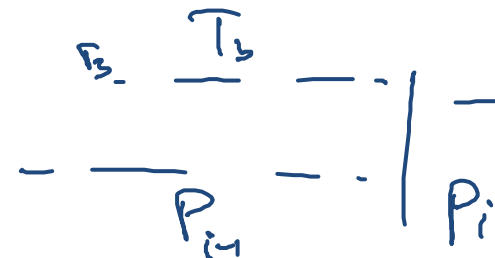
2.  $t_b$  is deleted:

$$E_{b,i} = E_{b-1,i} + c(t_b, \varepsilon)$$



3.  $p_i$  is inserted:

$$E_{b,i} = E_{b,i-1} + c(\varepsilon, p_i)$$



# Approximate String Matching

Recurrence relation (unit cost model):

$$E_{b,i} = \min \begin{cases} E_{b-1,i-1} + \mathbf{1} \\ E_{b-1,i} + \mathbf{1} \\ E_{b,i-1} + \mathbf{1} \end{cases} \quad \text{if } t_b = P_i$$

Base cases:

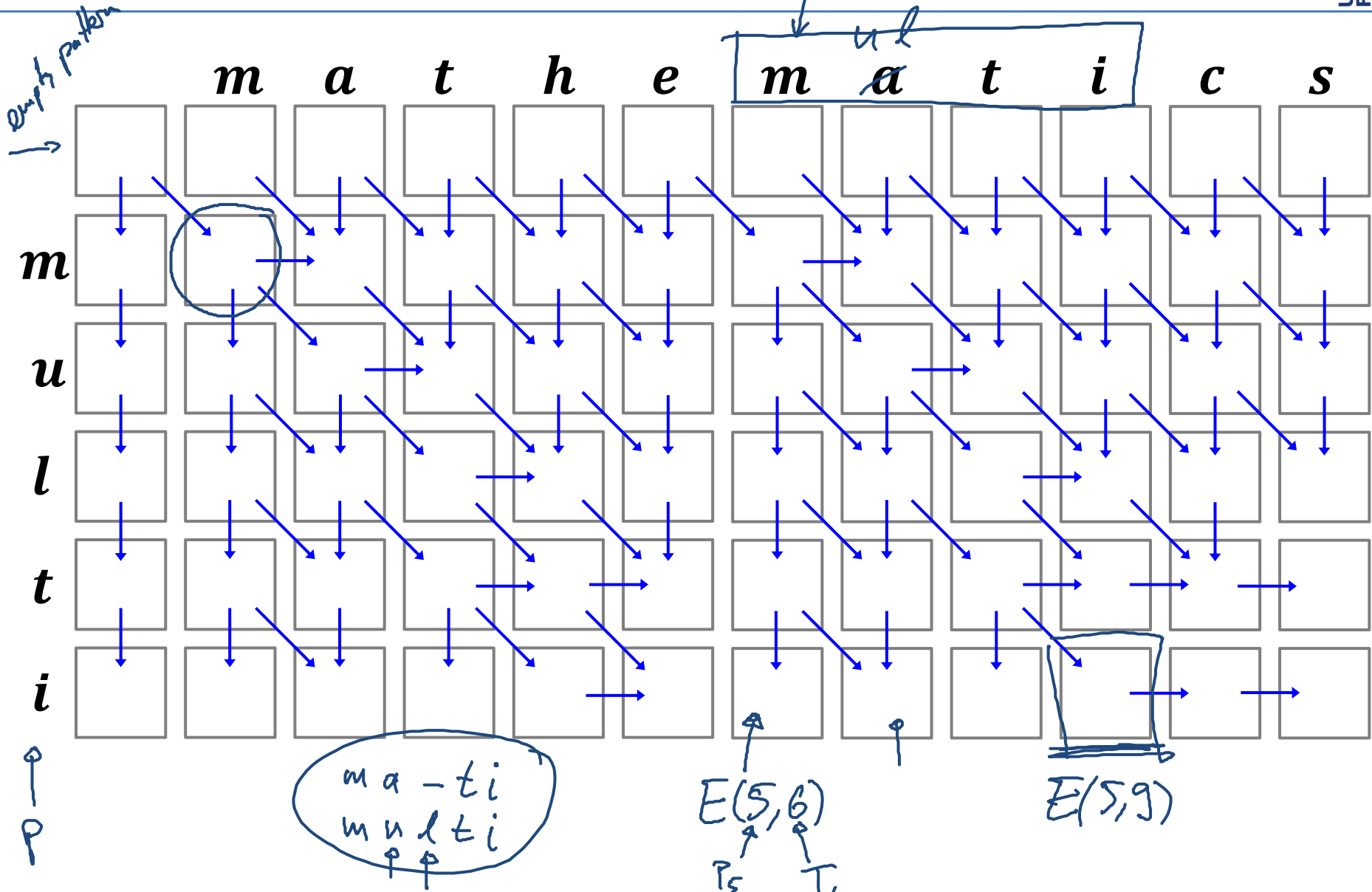
$$\begin{aligned} \rightarrow E_{0,0} &= 0 \\ E_{0,i} &= i \\ \rightarrow E_{i,0} &= 0 \end{aligned}$$

$$T = z, P_i$$

$$T_i = t_1 \dots t_i$$

$$P_0$$

# Example





# Approximate String Matching

---

- Optimal matching consists of optimal sub-matchings

$$m \ll n$$

- Optimal matching can be computed in  $O(mn)$  time
- Get matching(s):
  - Start from minimum entry/entries in bottom row
  - Follow path(s) to top row
- Algorithm to compute  $E(b, i)$  identical to edit distance algorithm, except for the initialization of  $E(b, 0)$

# Related Problems from Bioinformatics



## Sequence Alignment:

Find optimal alignment of two given DNA, RNA, or amino acid sequences.

```
G A - C G G A T T A G
G A T C G G A A T - G
```

## Global vs. Local Alignment:

- *Global alignment*: find optimal alignment of 2 sequences
- *Local alignment*: find optimal alignment of sequence 1 (pattern) with sub-sequence of sequence 2 (text)